

(D)

11 600

THE RESIDENCE OF THE PARTY OF T

Ju

RADC-TR-81-183, Vol IV (of four) Final Technical Report July 1981

A STATISTICAL THEORY OF COMPUTER PROGRAM TESTING

Polytechnic Institute of New York

Arthur E. Laemmei

APPROVED FOR PUBLIC RELEASE; DISTRIBUTION UNLIMITED

F HLE CORY

ROME AIR DEVELOPMENT CENTER
Air Force Systems Command
Griffiss Air Force Base, New York 13441

81 10 5 002

A

tes

This report has been reviewed by the RADC Public Affairs Office (PA) and is releasable to the National Technical Information Service (NTIS). At NTIS it will be releasable to the general public, including foreign nations.

RADC-TR-81-183, Vol IV (of four) has been reviewed and is approved for publication.

Court Links and the Court of th

APPROVED:

ROCCO F. IUORNO Project Engineer

APPROVED:

JOHN J. MARCINIAK, Colonel, USAF Chief, Information Sciences Division

FOR THE COMMANDER: Oslu G. Kuss

JOHN P. HUSS

STO 15 THE REPORT OF SERVICE PARTY.

Acting Chief, Plans Office

408704

If your address has changed or if you wish to be removed from the KAUL mailing list, or if the addressee is no longer employed by your organization, please notify RADC (ISIE) Griffiss AFB NY 13441. This will assist us in maintaining a current mailing list.

Do not return this copy. Retain or destroy.

Comment of the state of the sta

UNCLASSIFIED

SECURITY CLASSIF-CATION OF THIS PAGE When Data Entered)

REPORT DOCUMENTATION PA		READ INSTRUCTIONS BEFORE COMPLETING FORM
	SOUT ACCESSION NO	ARCIPTENT'S CATALOG NUMBER
RADC $\frac{1}{7}$ TR-81-183, Vol IV (of four)		C and a second
年、TITLE (and Subtitle)		5. TYPE OF REPORT & PERIOD COVE
SOFTWARE MODELING STUDIES . Vell	•	Final Technical Repert
A STATISTICAL THEORY OF COMPUTER PR	ROGRAM	January 78 - October 80
TESTING		S PERFORMING DIG REPORT NUMBER
T. AUTHORIS	11/1/PON	EE-80-004/
/ / / / / / / / / / / / / / / / / / /	مرد المسلم	SONTRACT OR GRANT NUMBER(S)
Arthur E. Laemmel	1:	F30602-78-C-0057
9 PERFORMING ORGANIZATION NAME AND ADDRESS		D PROGRAM ELEMENT PROJECT TA
Polytechnic Institute of New York		APERS WORK JUST NUMBERS
333 Jay Street		0.61102F
	11	2304.1401
Brooklyn NY 11201		'2 PEPORT CATE
		duly 1981 /
Rome Air Development Center (ISIE)	11	13 NUMBER OF PAGES
Griffss AFB NY 13441		36
14. NO. TORING AGENCY NAME & ADDRESS(If different from	n Cintrolling Office;	IS SECURITY CLASS for this report
		1
Same / Y		UNCLASSIFIED
$\mathcal{U}_{\mathcal{L}}$	e f	154. DECLASSIFICATION DOWNGRADIN
Approved for public release; distri		
		jN/A
Approved for public release; distri		jN/A
Approved for public release; distri TO DISTRIBUTION STATEMENT (of the abstract entered in B) Same 18. SUPPLEMENTARY NOTES) k 20 - if different in	jN/A
Approved for public release; distri T DISTRIBUTION STATEMENT (of the abstract entered in B) Same) k 20 - if different in	jN/A
Approved for public release; distri TO DISTRIBUTION STATEMENT (of the abstract entered in B) Same 18. Supplementary Notes RADC Project Engineer: Rocco F. In	orno (1811)	ited
Approved for public release; distri T DISTRIBUTION STATEMENT (of the abstract entered in B). Same 18. SUPPLEMENTARY NOTES RADC Project Engineer: Rocco F. In 9. KEY MORDS (Continue on reverse side of necessary and ide Program Testing Prog	orno (1811) ntily by block number	ited
Approved for public release; distri 17 DISTRIBUTION STATEMENT (of the abstract entered in B). Same 18. SUPPLEMENTARY NOTES RADC Project Engineer: Rocco F. In 9. KEY MORDS (Continue on reverse side of necessary and ide Program Testing Prog Testing Strategy Rand	orno (1511) ntily by block number ram Errors on Testing	ited om Report)
Approved for public release; distri T DISTRIBUTION STATEMENT (of the abstract entered in B). Same 18. Supplementary Notes RADC Project Engineer: Rocco F. In 9. KEY MORDS (Continue on reverse side if necessary and ide Program Testing Prog Testing Strategy Rand Probability of Failure Soft	orno (1811) ntily by block number	ited om Report)
Approved for public release; distri 17 DISTRIBUTION STATEMENT (of the abstract entered in B). Same 18. SUPPLEMENTARY NOTES RADC Project Engineer: Rocco F. In 9. KEY MORDS (Continue on reverse side of necessary and ide Program Testing Prog Testing Strategy Rand	orno (1511) ntily by block number ram Errors on Testing	ited om Report)
Approved for public release; distri 17 DISTRIBUTION STATEMENT (of the abstract entered in B). Same 18. Supplementary notes RADC Project Engineer: Rocco F. In 9. KEY MORDS (Continue on reverse side if necessary and ide Program Testing Prog Testing Strategy Rand Probability of Failure Soft Program Correctness	orno (1811) nully by block number ram Errors om Testing ware Error M	ited om Report)
Approved for public release; distri T DISTRIBUTION STATEMENT (of the abstract entered in B). Same 18. Supplementary Notes RADC Project Engineer: Rocco F. In 9. KEY MORDS (Continue on reverse side if necessary and ide Program Testing Prog Testing Strategy Rand Probability of Failure Soft	orno (1811) nutly by block number ram Errors on Testing ware Error Maring the rived to eating tested is those test of	odels eliability of computer mate the number of tests correct with a given

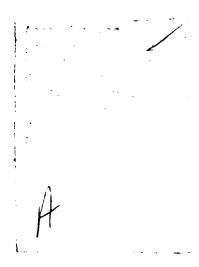
CONTRACTOR OF THIS PAGE When Date Entered)

408704

TABLE OF CONTENTS

Abstract

(). J	Introduction	1
2.0	Definitions and Objectives 2.1 Definitions 2.2 Objectives	3 3 3
3.0	Models 3.1 Elementary Model 3.2 Model with Unequal Probabilities 3.3 Model with Statistical Dependence 3.4 Illustrative Special Case 3.5 Optimum Testing Strategy (in Section 3.3) 3.6 Interpretation of "Bug" in the Last Model	4 4 6 8 13 13
4.0	Random Testing 4.1 Random Test Models 4.2 Optimum Testing Strategy (in Section 4.1)	17 17 22
5.0	Alternative Definitions of Error	23
6.0	Conclusions and Comments	24
7.0	References	25



LIST OF FIGURES

Figure 1	Plot of P_e vs. W from the Model of Equation 2	5
Figure 2	Plot of σ_{ji}	11
Figure 3	Bug Covered vs. Number of Tests	15
Figure 4	Probability of Error vs. Number of (Random) Tests	20
Figure 5	Specific Examples of Error Probability Curves	21

A STATISTICAL THEORY OF COMPUTER PROGRAM TESTING

Arthur E. Laemmel

Abstract

Most computer programs are tested with some of the possible sets of input data, but few can be tested with all possible input data. Passing such a partial test cannot insure that the program will always function correctly; we can perhaps say that the probability of failure is less than a specified amount. It is the purpose of this report to derive several formulas for the above probability. Also, in some cases an optimum testing strategy can be derived which minimizes that probability of failure.

1.0 Introduction

There are three methods for maximizing the reliability of computer programs: (1) use a systematic procedure which makes it difficult for errors to occur during the writing of the program, (2) prove that the program works correctly by some formal or automatic process, and (3) test and debug the program thoroughly before passing it on to the user. Most programmers will use some combination of these methods, and in fact some procedures involve elements of more than one method. The present report emphasizes testing, but first some remarks will be made about program writing and proving.

- (1) Writing Correct Programs: While no one intentionally inserts errors in his program, it is undoubtedly true that many people would produce more reliable programs with less effort if they were taught better programming techniques. However, it seems obvious that the average programmer should test his programs even if he exercises the maximum of care and uses the best techniques.
- Proving Program Correctness: There are several reasons why a formal procedure for proving program correctness cannot be relied on to insure absence of errors in practical situations: (i) a uniform algorithm for proving the correctness of an arbitrary program can be shown to be impossible, being essentially equivalent to Turing's halting problem; (ii) even for socyable sub-classes of the correctness-proving problem, the usual method (some improvement on Herbrand search) is so time consuming as to be impractical; (iii) there is always a possibility of error in the proving program, or in applying it to the program being tested.
- (3) Testing Computer Programs: In view of the difficulties of validating a computer program by programming techniques or formal proof methods, it is believed that some amount of testing will always be necessary. The purpose of this report is to describe a model which shows the relationship between errors of different types and the probability that they will cause a program to fail, and also to suggest optimum testing methods which minimize the probability of program failure.

Throughout this report it is assumed that a computer program can be tested and that it will either pass or fail the test. It is also assumed that a tester and a user will interpret failure of the program in exactly the same way. For example, failure might mean one of the following:

- i) the program "bombs" completely.
- ii) an obviously wrong answer is given.
- all a number at inaccurate in the least op to antideat

iv) the numbers are correct, but the format is wrong.

v) the program works perfectly, but a side effect causes failure in another program.

It must also be decided whether the program alone is being tested, or whether the program and algorithm is being tested. This report is directed mainly to the latter, but the results can be suitably interpreted so as to apply to the former.

The extent to which a computer program (possibly including the underlying algorithm) can be tested varies from application to application and is usually not a yes - no situation. Whether or not a statistical model applies to a particular case depends very strongly on the tester's knowledge of just what answer should be produced by the program. Some of the many possibilities of the user's prior knowledge of the correct answer are:

- 1. The exact answer is known beforehand. An example would be a math package for a new computer. Accurate tables of cos, arctan, log, etc. have been available for many years.
- 2. A proposed answer can be checked to see if it is a correct answer, but the correct answer is not known beforehand. Programs which calculate the roots of transcendental equations are examples of this category.
- 3. Answers for certain special input values are known beforehand. This is very common. For example, if a program is supposed to output the capacitance of an arbitrary two conductor transmission line, it would be natural to test it for coaxial cylinders.
- 4. The answer is not known beforehand, nor can it be accurately checked for any combination of input values. However, certain consistency relations among different outputs are known. For example, it may be obvious that the program should generate an output which is a monotonically increasing function of the input. A test which detects a decrease indicates an incorrect program, but no amount of testing can indicate a correct program.
- 5. Absolutely nothing is known about the answers which should be produced. This is probably very rare. Even here, the theory presented in this report applies if the program "bombs," i.e., a fatal or non-fatal run-time error message is produced.

The identification of an error is often not unique. This is illustrated by the fact that error messages from a compiler or run-time monitor can direct the user away from what he considers as the error. For example, if the compiler says "line 15, operator missing" the error might really be "line

14, statement terminator missing." In many cases the syntax can be corrected in several places to get the program thru the compile stage, but one place usually contains the reported error, and another place is usually where the error should be corrected. Similar comments can be made about semantic errors detected at run time. In some cases an error might equally well be corrected in one of several places. For example, a common PL/1 error might be corrected by either declaring K to be a FLOAT number, by using XK instead, or by forcing a necessary type conversion. This type of ambiguity in defining the error is not believed to cause any difficulties with the formulas developed in this report, provided the parameters are interpreted correctly.

2.0 Definitions and Objectives

2.1 Definitions. Some basic aspects of the testing process apply equally well to a computer program or to a physical device. For this reason the program or device being tested will sometimes be called the module. After the module has been constructed, it is checked by a tester and then employed the a user. The probability of error, $P_{\rm e}$, which occurs during use is given by

$$P_{e} = P_{m} \Gamma_{u} \tag{1}$$

where P_{m} is the probability that the tester misses all of the residual bugs in the module, and P_{tt} is the probability that the user then encounters one of the overlooked bugs. If exhaustive testing is possible then $P_{m}=0$, since it is assumed that if a bug is found it then is corrected and the whole process is started again. In most cases of interest exhaustive testing is not possible or practical. If there are no bugs then all of the probabilities are zero and this possibliity appears as a special case in the analysis to follow. Some alternate definitions of "probability of error" are given in Section

diectives. Perfore presenting the detailed mathema of results, the objectives of this work will be stated more explicitly. Roughly, the aims are to provide an estimate of the number of tests required to say that the program being tested is correct with a given probability, and to provide a strategy so that a given number of tests are chosen most effectively.

The first objective is met by deriving a functional relationship between the number of tests and the probability of error. There we several possibilities for defining probability of errors, the one used note might be called the "probability of embarrassment" for the tester is not embarrast at if

i) he discovers a bug and returns the program to the writer

he approves the program in spite of its having one or more bugs, but the user does not encounter one of the remaining bugs. The tester is embarrassed if

he approves the program and then the user encounters an undetected bug. If the model is designed appropriately, the probability of error decreases as the number of tests is increased according to Eq. 30 below.

The second objective is met by choosing those tests which minimize the expression derived for probability of error (defined as above) while keeping the number of tests fixed. Specifically, the principal question which must be answered here is: should testing be concentrated on input data combinations which are most likely to be chosen by the user, or input data combinations with a large a priori probability of causing program failure? The answer is that tests should be applied to both of these combinations of input data in a ratio which can be calculated from Eq. 33 below.

Actually, several simpler models are also analyzed before that described by Eqs. 30 and 33. All of these expressions contain many parameters, and curves are plotted for selected values. No real data were used for the parameters, but this should certainly be done in the future.

3.0 Models

3.1 Elementary Model A simple case might be the following: The module has N possible input values, and each of these are equally likely to be chosen by the tester or by the user. Of these input values, W cause improper functioning of the module, but neither the tester nor the user knows which inputs cause errors or even how large W is. The tester chooses t inputs at random without keeping a record of inputs previously tested, i.e., sampling with replacement. Under these circumstances $P_u = W/N$, $P_m = (1 - W/N)^t$ and

$$P_{e} = \frac{W}{N} \left(1 - \frac{W}{N} \right)^{t} \le \frac{W}{N} e^{-Wt/N}$$
 (2)

A plot of P_e vs. W is displayed in Fig. 1. If the testing is to do any good, i.e., to reduce P_e significantly less than W/N, then it is necessary that

In many applications it is found that satisfying the inequality of Eq. (3) requires a very large number of tests t, and that our intuitive feeling is that $P_{\underline{e}}$ is acceptably small in spite of testing using far fewer tests than indicated.

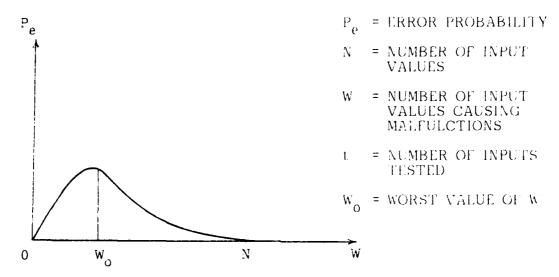


FIGURE 1. PLOT OF P_{e} VS. W FROM THE MODEL OF EQUATION 2.

It will continue to be assumed that the tester passes no information to the user about which inputs were tested. It will also be assumed that each test either succeeds or fails and that there is no additional information which would permit sequential sampling methods. Sampling without replacement would slightly lower $P_{\mathbf{m}}$ to

$$P_{m} = \begin{cases} \frac{(N-W)!(N-t)!}{(N-W-t)!} & t < N-W \\ 0 & t > N-W \end{cases}$$
 (4)

Note that this method requires that the tester keep a receid of inputs already tested, or that he avoids duplication by other means.

The particular type of error to which $P_{\tilde{e}}$ pertains must be borne in mind to avoid confusion with other possible definitions of error. If W=N all inputs to the module cause malfunctions but $P_{\tilde{e}}=0$ according to Equations (2) and (4). This is true because the tester removes user bugs with each test, and continues until all tests are exhausted. In this case the tester always rejects the module (provided only that t > 0) and so the user cannot experience a malfunction. Note that $P_{\tilde{e}}$ is neither the probability (user has a malfunction) nor the probability (user has a malfunction) fescer accepts module. Rather, $P_{\tilde{e}}$ is the probability tuser has a malfunction indicated and tester accepts module. If the number of tests is fixed at the and if testing with the contribution of the standard malfunction and tester accepts module. If the number of tests is fixed at the and if testing with the contribution is done, then the value of Washard measures $P_{\tilde{e}}$ is

$$W_{O} = \frac{N}{t+1} \tag{5}$$

From Eq. (5) or Fig. 1, it can be seen that as W \rightarrow 0 then P \rightarrow 0 also. This is so because for small W there is a small chance that the user will encounter a faulty input. Similarly as W \rightarrow N then P \rightarrow 0 also because there is little chance that the tester will accept the module. Of course, the last case is undesirable for reasons other than the value of Pe, i.e., the user has a small probability of receiving a released program.

- 3.2. Model with Unequal Probabilities The model described in the preceding section is too simple to apply to most practical testing situations: some inputs are more likely to fail than others, and the probability of one input failing may not be independent of another input failing. Often a single bug may cause many inputs to fail. The tester may not choose the inputs to be tested randomly, but rather in such a way as to utilize his knowledge of the prior failure probabilities. The user may not be free to choose more reliable inputs; in fact, he may be constrained by the problem to use less reliable inputs. The model to be described here includes three events, the last two being independent of each other but dependent on the first.
 - (1) A programmer constructs a module which has an error pattern α with probability $\underline{P}(\alpha)$. α might be a binary vector $(\alpha_1, \alpha_2, \ldots, \alpha_N)$ with $\alpha_i = 1$ meaning input i fails and $\alpha_i = 0$ meaning input i functions correctly.
 - (2) A tester tries certain inputs to the module and accepts the module with probability $R(\text{accept} \mid \alpha)$. The tester passes no information concerning which inputs were tested to the user.
 - (3) A user selects one of the inputs and the module fails with probability $Q(\text{fail} \mid \alpha)$. The error probability defined previously is now given by

$$P_{e} = \sum_{\alpha \in A} P(\alpha) \ \underline{Q}(\text{fail} \mid \alpha) \ \underline{R}(\text{accept} \mid \alpha)$$
 (6)

where A is the set of all possible error patterns.

Fo illustrate, if there are N input values then A consists of $2^{\mathbb{N}}$ elements. Assume P(α) is 0 for all A except for the first w:

$$\alpha_{W} = (1, 1, \dots, 1, 0, 0, \dots, 0)$$
 $(6a \cdot W) \rightarrow (N - W) \rightarrow (6a \cdot W)$

and let the probability of the user selecting input i be \mathbf{q}_i . Then

$$\underline{Q}(\text{fail} \mid \alpha_{W}) = \sum_{i=1}^{W} q_{i}$$

Assume the tester selects his inputs randomly, choosing input i with probability $\mathbf{r_i}.^{*}$. Then

$$\underline{R}(\text{accept } \mid \alpha_W) = \lim_{j=1}^{W} (1 - r_j)$$
 (7)

$$P_e = (\sum_{i=1}^{W} q_i) \prod_{j=1}^{W} (1 - r_j)$$

If $q_i = 1/N$ and $r_j = t/N$ this reduces (approximately) to the elementary model given above:

$$P_{e} = \frac{W}{N} \left(1 - \frac{t}{N}\right)^{W} \tag{8}$$

Another, more useful, form for $P(\alpha)$ than Eq. (6a) is obtained by assuming that the i^{th} input malfunctions with probability p_i . Then

$$\underline{P}(\alpha) = \prod_{i=1}^{N} p_{i}^{\alpha_{i}} (1 - p_{i}^{1-\alpha_{i}})$$

$$Q(\text{fail} \mid \alpha) = \sum_{j=1}^{N} \alpha_{j} q_{j}$$

$$R(\mathbf{a} \mid \alpha) = \prod_{i=1}^{N} (1 - r_i)^{\alpha_i}$$

The formula also applies if r_i is given the interpretation "input i is tested and the response is noted to be wrong by the tester." Specifically $r_i = 0$ might mean either that input i was not tested, or that it was tested and an error was not detected.

The two products can be combined in evaluating P_{e} :

$$P_{e} = \sum_{\alpha_{1}} \sum_{\alpha_{2}} \dots \sum_{\alpha_{N}} \sum_{j=1}^{N} \alpha_{j} q_{j} \prod_{i=1}^{N} F_{i}(\alpha_{i})$$

where

$$\Gamma_i(\alpha_i) = p_i^{\alpha_i} (1-p_i)^{1-\alpha_i} (1-r_i)^{\alpha_i}$$

This reduces to

$$P_{e} = \prod_{j=1}^{N} (1 - p_{j} r_{j}) \sum_{i=1}^{N} q_{i} \frac{p_{i} (1 - r_{i})}{1 - p_{i} r_{i}}$$
(9)

If the tester selects t inputs deterministically, and if the inputs are permuted so that these occur first, then

$$P_{e} = \prod_{i=1}^{t} (1-p_{i}) \sum_{i=t+1}^{N} p_{i}q_{i}$$
 (10)

An optimum testing strategy is obtained if the inputs are permuted so that the expression is minimized. The first factor suggests testing inputs with the largest p_i , but the second factor suggests testing inputs with the largest $p_i q_i$. Let the above expression be abbreviated as $P_e = P_m P_u$ and note that P_u is the probability of the user getting an error on an untested input. Consider the effect of adding one more input to the test.

Let

$$P'_{e} = P_{m}(1-p_{k})(P_{u}-p_{k}q_{k})$$

 $P'_{e} \approx P_{e} (1-\frac{(q_{k}+P_{u})p_{k}}{P_{u}})$

Thus, the criterion is to select the input with the largest

$$(q_k + P_u)p_k \tag{11}$$

As can be seen, this provides a weighted compromise between selection on the basis of ρ_k or $q_k \rho_k$ alone.

3.3 Model with Statistical Dependence Computer programs usually fail for a whole set of input values as a result of a single bug. It is more

realistic to assume that failures due to different bugs are statistically independent rather than failures for different input values. For example, a single oversight might cause a square root program to fail for all negative numbers. Let

$$\sigma_{ij} = \begin{cases} 1 \text{ if the } j^{th} \text{ bug causes failure for input } i \\ 0 \text{ if the } j^{th} \text{ bug doesn't affect input } i \end{cases}$$

If β_i is the probability of the j^{th} bug, if M is the number of possible bugs, and if θ_j (j = 1,2,...,M) is the pattern of actual bugs, then (note definitions of P, Q, R in the probability of event statements (1), (2), and (3) of Section 3.2)

$$P_{e} = \sum_{\theta} \underline{P}(\theta) \ \underline{Q}(f \mid \theta) \ \underline{R}(a \mid \theta)$$
 (12)

This is analogous to the corresponding formula in α given above. Here,

$$\underline{P}(\theta) = \prod_{i=1}^{M} \beta_{i}^{\theta_{i}} (1 - \beta_{i})^{1 - \theta_{i}}$$

$$\underline{Q}(f \mid \theta) = \sum_{j=1}^{N} q_{i} \{1 - \prod_{k=1}^{M} (1 - \sigma_{jk} \theta_{k})\}$$

$$\underline{R}(\alpha \mid \theta) = \prod_{i=1}^{M} \prod_{\ell=1}^{N} (1 - r_{\ell})^{\sigma_{\ell} \ell} \theta_{i}$$
(13)

Combining and rearranging gives

$$P_{e} = \sum_{j=1}^{N} q_{j} \sum_{\theta} A_{j}(\theta_{1}, \theta_{2}, \dots, \theta_{M}) \prod_{i=1}^{M} F_{i}(\theta_{i})$$

where

$$\Gamma_{i}(\theta_{i}) = \beta_{i}^{\theta_{i}} (1-\beta_{i})^{1-\theta_{i}} \prod_{\ell=1}^{N} (1-r_{\ell})^{\sigma_{\ell}i\theta_{i}}$$

and

$$A_j(\theta_1,\ldots,\theta_M) = 1 - \prod_{k=1}^{M} (1-\sigma_{jk}\theta_k)$$

The summations over θ_i are for only two values (0,1) and can be carried out to give

If deterministic testing is used (r=0,1) over the first t inputs:

$$P_{e} = \prod_{i=1}^{M} (1-\beta_{i}) \sum_{j=t+1}^{N} q_{j} \left[1 - \prod_{i=1}^{M} (1-\sigma_{ji}\beta_{i})\right]$$

Where Π' refers to a product over terms involving a value of i such that σ_{ij} = 1 for at least one value of j in the range j = 1,2,...,t, and where Π'' refers to the other values of i. The bugs can be permuted so that $1 \leq i \leq \tau$ implies σ_{ji} = 1 for at least one value of j from 1 to t and $\tau < i \leq M$ implies σ_{ji} = 0 for all values of j from 1 to t. Then

$$M$$
 τ M M Π' = Π and Π'' = Π $i=1$ $i=1$ $i=\tau+1$

The problem is then to minimize

$$P_{e} = \prod_{i=1}^{\tau} (1-\beta_{i}) \sum_{j=t+1}^{N} q_{j} \left[1 - \prod_{i=\tau+1}^{M} (1-\sigma_{ji}\beta_{i})\right]$$
 (15)

A graphical interpretation of the above is portrayed in Fig. 2, and illustrates how different inputs excite various bugs. For example, input 2 excites bugs numbered 6,7,8 and 9. Bug numbered 8 can only be discovered through the application of inputs 2 or 4.

A clearer picture of how P_e depends on the amount of testing, t, can be obtained by deriving an upper bound from Eq. (15). For most cases of interest this upper bound will be tight, i.e., a good approximation. Multiply thru by the first product in Eq. (15):

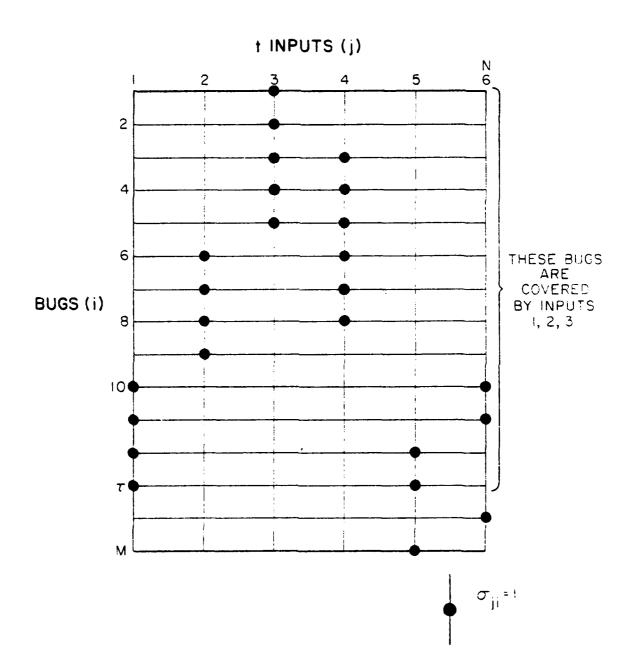


FIGURE 2. PLOT OF σ_{ji}

$$P_{e} = \sum_{j=t+1}^{N} q_{j} \left\{ \prod_{i=1}^{\tau} (1-\beta_{i}) - \prod_{i=1}^{M} (1-\gamma_{ji}\beta_{i}) \right\}$$
 (16)

where

$$y_{ji} = \begin{cases} 1 & 1 \le i \le \tau \\ \sigma_{ji} & \tau < i \le M \end{cases}$$

Now observe the following:

$$-\prod_{i=1}^{M} (1-\gamma_{ji}\beta_i) = -\prod_{i=1}^{M} (1-\beta_i) \quad \text{since} \quad \gamma_{ji} \leq 1$$

$$\sum_{j=t+1}^{N} q_j - Q(t) \le 1 \quad \text{(note } Q(t) = 1 - \frac{t}{N} \stackrel{\sim}{\sim} 1 \text{ if } q_j = \frac{1}{N} \text{)}$$

$$\begin{array}{ccc} & & \tau & \\ \tau & & i=1 \\ \Pi & (1-\beta_i) \leq e \end{array}$$

Combining these:

$$P_{e} \leq \bar{P}_{e} = e \qquad -R_{o}$$

$$(17)$$

where

$$R_{0} = \prod_{i=1}^{M} (1-\beta_{i}) \approx e$$

From the way the σ_{ji} matrix was permuted, it can be seen that t \rightarrow M (monotonically), and from Eq. (16) these, in turn, imply $P_e \rightarrow 0$ (monotonically).

- If Eq. (17) is a good approximation to Eq. (16), then the rate at which $P_{\rm e} \neq 0$ can be seen to be exponential.
- 3.4. Illustrative Special Case. A simple special case will illustrate the above formulas, and is useful in getting a rough idea of the number of tests and the probability of error. Suppose that each of the M bugs occurs with the same probability of error, β , and affects the same number of input value, b. Suppose further, that the pattern of input errors is the most difficult to detect with the given number of tests, t, i.e., that the error subsets are as "disjoint" as possible. If the tests are distributed most effectively over the inputs, each test will detect Mb/N bugs. The testing will result in M Mbt/N undetected bugs, and $\tau = \text{Mbt/N}$ (assuming bt \leq N). Thus

$$P_{\epsilon} = \frac{\beta M D t}{N} = -\beta M \tag{18}$$

- As $t \sim N/h$ this move $T_{\bf e} = 0$. It order to make $T_{\bf e}$ small comparable to New More general assumptions concerning the parameters, and with a publishmode of tests, t might be a much smaller fraction of N.
- 3.5 Optimum Testing Strategy (in Section 3.3) For a given t, that selection of inputs to be tested which minimizes P_e will be defined as the optimum testing strategy. From Eq. (15), and intuition, one might be tempted to test the t inputs most likely to be chosen by the user, thus minimizing the factor Σq . However, for most cases of interest, this approach would be futile. If the q are even very roughly equal, the number of possible inputs N is so hage that it would be impossible to test enough of them to reduce

significantly below unity. Essentially, the $q_{\hat{j}}$ only influence the testing strategy thru the fact that the square bracketted term in Eq. (15) depends on i.

The first term in Eq. (15) is the most important for values of the parameters of most interest, and as was shown above, it is minimized (approximately) by choosing test values to maximize

$$E = \sum_{i=1}^{\tau} \beta_i$$
 (19a)

$$E = \tau \beta$$
 if $\beta_i = \beta$ (19b)

If the β_i are equal, P_e is then minimized (approximately) by maximizing $\tau.$ This is essentially the covering problem of switching theory and operations research.
[1,2] A set of inputs i ϵ I is said to cover a set of bugs j ϵ J if σ_{ji} = 1 for every j in J for at least one i in I. Referring to Fig. 2, input 4 covers bugs 3, 4, 5, 6, 7 and 8; and input pair 1,4 covers bugs 3, 4, 5, 6, 7, 8, 10, 11, 12 and 13. These inputs are optimum in that they cover the most bugs possible. Thus, using these optimum choices, τ = 6 and 10 for t = 1 and 2. However, the optimum set of 3 inputs is not obtained by adding that input which would cover the most additional bugs. Inputs 1, 3, 4 cover only 12 bugs, but inputs 1, 2, 3 cover 13 bugs. The covering problem is usually stated as minimizing the number of tests required to cover all bugs (in the present terminology), and no general algorithm for its solution is known. The process described above, repeatedly adding that test which causes the greatest increment in bugs covered, is usually referred to as the "heaviest-first algorithm."

A sketch of τ vs. t for the present example is shown in Fig. 3.

Note that the number of covered bugs added at each step is a decreasing function of t. If t is very large, the bugs being covered might be those affecting only a single input, such as divide-by-zero bugs.

The rate of increase of P_e with increasing t is certainly of great interest, and one simple functional dependence can be given using the above ideas. Define τ_1 , τ_2 ,... as those values corresponding to t=1,2,... when optimum or near-optimum testing strategy is used. Eq. (19a) can be rewritten

$$E_{t} = \sum_{i=1}^{\tau_{t}} \beta_{i} = \sum_{k=1}^{\tau} B_{K}$$
 (20)

where

$$B_K = \sum_{\tau_{k-1}+1}^{\tau_k} \beta_i, \quad \tau_0 = 0$$

From Eq. (15) with the second and third terms approximated by unity, or from Eq. (17) with $R_{\rm o}$ approximated by zero:

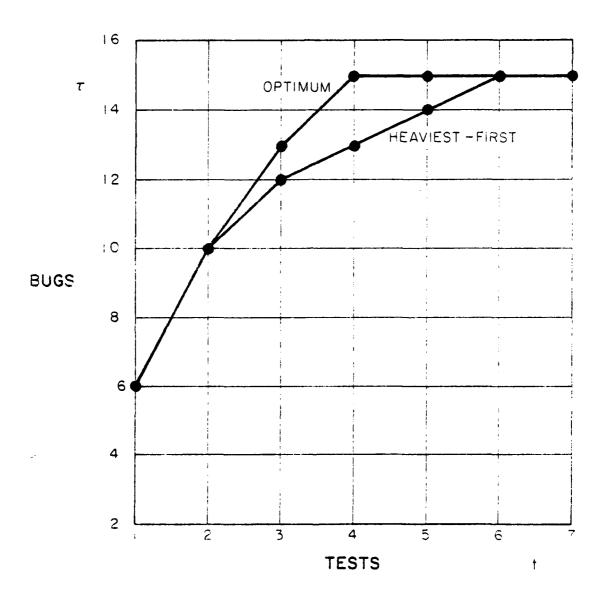


FIGURE 3. BUGS COVERED VS. NUMBER OF TESTS

$$P_{e}(t) \approx e^{-\frac{t}{2}BK}$$
(21)

Now if B_k is assumed to have a form which is a decreasing function of k and which can be finitely summed in closed form, a convenient relation can be found. For example, if $B_k = c/k$

$$\frac{t}{\sum_{k=1}^{\infty} \frac{c}{K}} \approx 0.5772c + c \ln t$$

$$P_{e}(t) \approx e^{-(c-\sqrt{t}/2c)}$$
(22)

This shows what might be expected in practice, a gradual decrease of P_e over many decades as t is increased over several decades. When some relations are plotted, eg. Eq. (2) or Eq. (18), they exhibit very sudden and deep drops in P_e when certain values of t are approached, and this is not the behavior to be expected except in the simplest programs. Of course, Eq. (22) does not have enough parameters, but similar formulas can be found which do have enough parameters.

- 3.6 Interpretation of a Program Bug There are many possible ways to interpret a bug.
- 1. Formally, a bug is simply a subset of input values which fail together when the program is used, e.g., the region ${\rm B}^2$ 4 AC < 0 in solving a quadratic.
- 2. A bug might be a careless typing error, such as A+B instead of A-B.
 - 3. A bug might be an incorrect statement or a defective subroutine.
- 4. A more flexible definition of bug should allow such things as omitting a check for dividing by zero. These bugs of omission are harder to handle with fixed M and matrix σ_{ii} .
- 5. A bug might be defined as a distinct path thru a flowchart. A popular testing procedure is to run at least once thru every such path, and if these are the only bugs included, this constitutes a complete cover and P_e = 0. Unfortunately, the model would then not be very realistic because P_e would never actually vanish except in the most trivial programs.

6. If programming can be identified with decisions, a bug is a wrong decision.

4.0 Random Testing

4.1 Random Test Models Some of the difficulties involved in relating the number of bugs discovered and the number of tests made can be avoided if the tests are chosen randomly rather than deterministically. Optimum strategies can still be found if the probabilities of choosing the inputs to be tested are not equal. The probability of error P_e for random testing has already been found in Eq. (14) above. Let \mathbf{s}_{ℓ} be the probability that the tester chooses the ℓ th input at any particular test. If the number of tests is t, then the probability that the ℓ th input is tested during the series of t tests is

$$r_{\varrho} = 1 - (1 - s_{\varrho})^{t}$$
 (23)

Note that the \mathbf{s}_{ℓ} sum to unity, but that each \mathbf{r}_{ℓ} can range from zero to one. It is helpful to define

$$A_{i} = \prod_{\varrho=1}^{N} (1-s_{\varrho})^{\sigma_{\varrho i}}$$
 (24)

which is the approximate probability of a single test missing the i'th bug. Note that A_i is certainly no greater than unity. Eq. (14) can now be rewritten in terms of s, A and t.

$$P_{e}(t) = \prod_{i=1}^{M} (1 - \beta_{i} + \beta_{i} A_{i}^{t}) \sum_{j=1}^{N} q_{j} \left[1 - \prod_{i=1}^{M} (1 - \frac{\sigma_{ji} \beta_{i} A_{i}^{t}}{1 - \beta_{i} + \beta_{i} A_{i}^{t}}) \right]$$
(25)

As the number of tests is increased the probability of error approaches zero

$$\lim_{t\to\infty} P_e(t) = 0$$

If no tests are made

$$P_{e}(0) = 1 - \sum_{i=1}^{N} q_{i} \prod_{j=1}^{M} (1 - \sigma_{ji} \beta_{i})$$
 (26)

It is easier to see the importance of different terms in Eq. (25) by examining the asymptotic form for large t:

$$\bar{P}_{e}(t) = \prod_{h=1}^{M} (1 - \beta_{h}) \sum_{i=1}^{M} \frac{\beta_{i}}{1 - \beta_{i}} Q_{i} A_{i}^{t}$$
(27)

where

$$Q_{i} = \sum_{j=1}^{N} q_{j} \sigma_{ji}$$

The quantity Q_i is the probability that the user encounters an error from the i'th bug. Since each s_ℓ is very small an exponential can be used to express A:

$$A_{\hat{i}} = \prod_{\ell=1}^{N} (1 - s_{\ell})^{\frac{\sigma_{\ell} \hat{i}}{\ell}} = e^{-S_{\hat{i}}}$$
(28)

where

$$S_{i} = \psi \sum_{\ell=1}^{N} s_{\ell} \sigma_{\ell i} \quad 1 \le \psi \le \frac{1}{1 - s_{\text{max}}}$$
 (29)

Note that if the tester uses the same probabilities as the user then $S_i = \psi Q_i$. If this assumption is made, if the β are small, and if Q_i takes on only a few distinct values Q_1, Q_2, \ldots, Q_n , then the asymptotic form of Eq. (17) can be approximated by

$$\bar{P}_{e}(t) \stackrel{\sim}{\sim} \sum_{k=1}^{n} \beta_{k} N_{k} Q_{k} e^{-Q_{k} t}$$
(30)

where

 $\rm N_{k}$ is the number of bugs with a probability of discovery $\rm Q_{k}$ (These bugs which have probility $\rm Q_{k}$ will be called a "block".)

and

 $\beta_{\mbox{$k$}}$ is the average probability of bugs in the k'th $\underline{\mbox{block}}$ being introduced by the program writer.

No loss of generality is obtained if the Q_k are arranged in order of decreasing magnitude. Eq. (30) clearly represents a decreasing function of t, but it is not so apparent that the decrease is more than 1/t than exponential over the range of interest. To see this, note that each term in Eq. (30) reaches a maximum as a function of Q_k at $Q_k = 1/t$. For each

t a certain term in the sum will be dominant and values of k both larger and smaller than this dominant term will contribute relatively small amounts. This leads to an approximate formula

$$\bar{P}_{e}(t_{k}) \approx \frac{\beta_{k} N_{k}}{t_{k} e} \tag{31}$$

where

$$t_k = \frac{1}{Q_k}$$

If the product $\beta_k N_k$ is approximately independent of k, then P_e will decrease approximately as $1/t_k$. The decrease is exponential below $t=1/Q_{min}$ but this value could be very large indeed. For example, if the input is a 32 bit number, and if a bug causes an error for a single input value, then $Q_{min} = 2^{32}$.

Such a bug would be very hard to detect by random testing, but by the same token it would be very unlikely to cause a user error unless the user favors that value for some reason. The behavior of Eqs. (30) and (31) is sketched in Fig. 4 and Fig. 5.

Fig. 5 shows the way in which the probability of error P_e decreases with the number of tests t. The parameters were chosen to illustrate the following cases:

- i) P_{ρ} decreases slowly at first, then rapidly
- ii) P_{ρ} decreases gradually throughout the range
- iii) Pe decreases rapidly with the first few tests, then many more tests are required to obtain a further decrease.

As might be expected, in case (i) above most bugs are difficult to detect, in case (ii) bugs are evenly distributed, and in case (iii) bugs are either very easy or very difficult to detect.

Program bugs have been divided into four, and in some cases five, classes, indicated by k in the following

 \boldsymbol{Q}_k is the probability of detecting a class k bug

 $\boldsymbol{\beta}_k$ is the probability that a user has introduced a class k bug

 $\mathbf{N}_{\mathbf{k}}$ is the number of possible class k bugs.

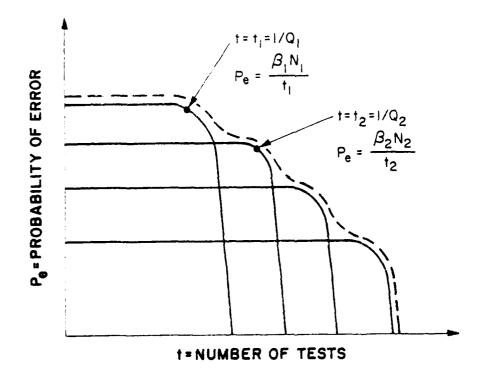


FIGURE 4. PROBABILITY OF ERROR VS. NUMBER OF (RANDOM) TESTS

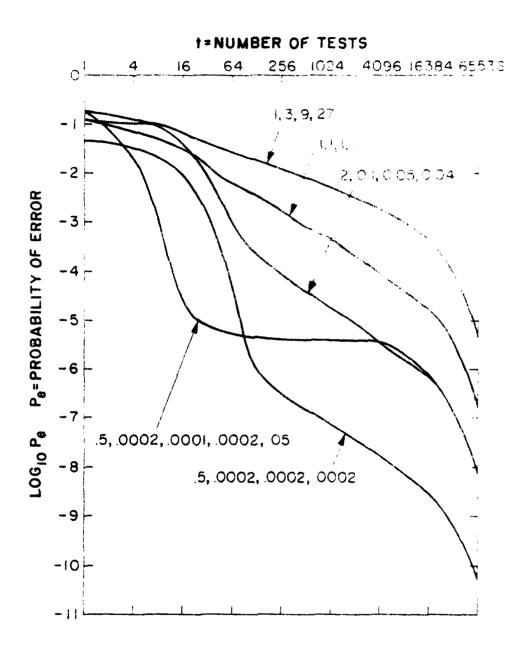


FIGURE 5. SPECIFIC EXAMPLES OF ERROR PROBABILITY CURVES (from Eq. (30), see text for explanation of parameters)

The parameter Q_k is always 0.1, .01, .001 and .0001, except where 5 classes are present in which case Q_k is 0.7, 0.1, .01, .01, .001, and .0001. Shown on the curves is the product $\beta_k N_k$ for each of the four or five classes. For example, the curve marked 1,1,1,1 applies to these 4 cases:

φ_k	$^{\beta}k$	^N k	$Q_{\mathbf{k}}$	$\beta_{\mathbf{k}}$	N k
0.1 .01 .001 .000	.05 .05	20 20 20 20 20	0.1 .01 .001 .0001	.01 .005 .002 .001	100 200 500 100
$Q_{\vec{k}}$	$\mathfrak{b}_{\mathbf{k}}$	^N k	$Q_{\mathbf{k}}$	$\beta_{\mathbf{k}}$	N _k _
0.1 .01 .00i .0001	.01 .001 .001	100 1000 1000 1000	0.1 .01 .001 .0001	.002 .01 .01 .002	500 100 100 500

Ease of detection, as reflected in $\mathsf{Q}_k\text{,}$ depends largely on the number of inputs a bug effects.

4.2 Optimum Testing Strategy (in Section 4.1). Even though the inputs to be tested are to be selected randomly, an optimum strategy exists in that the probabilities of testing various inputs can be chosen so as to minimize P_{e} . The asymptotic form expressed by Eq. (27) can be rewritten

$$\bar{P}_{e}(t) = (1-\beta_{n}) \sum_{i=1}^{M} \frac{\beta_{i}}{1-\beta_{i}} Q_{i}^{-ts_{i}}$$
 (32)

The quantities β_i and Q_i are fixed by the problem, but the s_i can be chosen to minimize \overline{P}_e , at least insofar as they depend on the s_ℓ (see Eq. (29)).

The optimum s_{ℓ} can easily be found for the special case where σ_{ji} = δ_{ji} , i.e., where each bug affects only one input. The problem is to minimize the asymptotic form (with some additional approximations):

$$f = \sum_{i+1}^{N} \beta_i q_i e^{-s_i t}$$

while keeping constant

$$1 = g = \sum_{i=1}^{N} s_i$$

The method of Lagrange multipliers gives the equation

$$0 = \frac{\partial f}{\partial s_k} + \lambda \frac{\partial g}{\partial s_k} = -\beta_k q_k t e^{-s_k t + \lambda}$$

After eliminating λ to insure that the s_i sum to unity, this gives

$$\mathbf{s}_{\mathbf{k}} = \frac{1}{N} + \frac{1}{t} \ln \left(c \mathbf{\beta}_{\mathbf{k}} \mathbf{q}_{\mathbf{k}} \right) \tag{33}$$

where

$$\ell n \ c = -\frac{1}{N} \sum_{i=1}^{N} \ell n \ \beta_i q_i$$

This solution shows that if a small number of tests are made, the tester should favor the inputs most likely to be used and most likely to be in error, but that if a large number of tests are to be made the tester should select the inputs with an essentially uniform distribution. Even where a non-uniform distribution of tests is indicated by Eq. (33), the tester should distribute his tests more uniformly than the user because of the logarithm (assuming the $\beta_{\bf k}$ do not deviate too markedly from uniformity).

The optimum testing strategy when σ_{ji} is a delta functon (each bug affecting only one input) can be quite misleading for the more general case. It may be quite impossible to make the s_i in Eq. (32) either independent of i, or (even logarithmically) proportional to the Q_i . Also, if one returns to the more general Eq. (25) or Eq. (17) and attempts to choose the s_ℓ to minimize P_e or \bar{P}_e , it is found that the best value for s_ℓ is either zero or one and that these values depend on an unrealistically detailed knowledge of $\sigma_{\ell i}$. This is really the same as the optimum deterministic testing discussed in Section 3.5.

5.0. Alternative Definitions of Error

The definition of error which was used above may not be suitable for all purposes, but related probabilities can easily be calculated from the equations given. Define two events as follows:

 $E_{\mathbf{u}}$ is the event that a user employs the module once and encounters a failure.

 E_{m} is the event that a tester operates the module t times and does not encounter a failure, i.e., the tester accepts the module. The subscript m is mnemonic for "missed," since bugs are always assumed to be present with some probability, and therefore $\{e_{m}\}$ is the probability that the tester has missed all bugs, i.e., incorrectly accepted the module.

Three different quantities which might be interpreted as the probability of user error are tabulated below (Table I) for the various models which have been analyzed. The effectiveness of the testing process can be judged by comparing the first with the last two. Note that the last, $\Pr\{E_{ij}|E_{m}\}$ is always greater than the second, $\Pr\{E_{ij}E_{m}\}=P_{e}$. This fact might lead to a model of conditional probability, $\Pr\{E_{ij}E_{m}\}$, in some cases where P_{e} is small due to tester rejection and where there is a large probability of module bug.

One is tempted to regard the goal here as minimization of $\Pr\{E_{ij}\}$, regarding the avoidance of user error as of most importance to the user. But remember this report seeks an optimum testing method, and that this is different from and does not preclude previously using an optimum programming method. The latter will minimize user error, but will not generally remove the need for testing in addition.

6.0 Conclusions and Comments

It is believed that defining what is meant by the probability of a program error, and presenting a model which permits its exact calculation (Eq. (15)) will provide the nucleus around which a theory of software reliability can be built. The purpose is not merely to get a formula into which numbers can be plugged to give a probability of error -- this does nothing to reduce errors. Rather it is anticipated that by classifying different types of bugs, errors and test results, and by showing how they interact programming systems can be improved and optimum testing procedures can be found.

Most reports on program testing select the test data so as to traverse each path in the program at least once. [3,4] The approach described in this report might seem to be directed at selecting test data according to the problem specification. A near-optimum strategy may be to select test data randomly, since this seems to be a fairly efficient solution to the covering problem, but to also make sure some test points are in various distinct data domains (e.g., B^2 - 4AC < 0) and to also make sure each program path is traversed. In addition, it is desirable to concentrate test data at points known to be more likely to be selected by the user. These requirements are often contradictory, but optimum compromises are suggested by Eq. (11), Section 3.5, and Section 4.2 above (provided the stated assump-

tions are met and the required parameters are available). Also formulas derived give some idea of how many tests need to be made to achieve a given reliability.

If further work is to be done along the lines of this report it is suggested:

- 1) Data be gathered so as to verify the derived relations between the number of tests and the error probability, especially as in Section 4.0.
- 2) Another model can be developed in which a partial knowledge of the bug-input matrix σ is assumed. This would lead both to a more easily applied optimum testing strategy, and to a method of sequential testing in which the last test points are chosen according to the results of the first test.

7.0 References

- 1) M.L. Balinski, "Integer Programming: Methods, uses, computations," Management Sci., Vol. 12, p. 253, November 1965.
- 2. Min-Wen Du, "A way to find a lower bound for a minimal solution of the covering problem," IEEE Trans. on Computers, Vol. 21, p. 317, 1972.
- 3. J.C. Huang, "An approach to program testing," Comput. Surv., Vol. 7, p. 113, September 1975.
- 4. W. Miller and D.L. Spooner, "Automatic generation of floating-point test data," IEEE Trans. on Software Engr., Vol. SE2, p. 223, September 1976.

Table I summarizes the models discussed.

TABLE E. Summary of the statistical models.

USer Count Course		ar.		[(\d\ \frac{1}{2} - 1 \cdot\ \frac{1}{2} - 1	$\frac{1}{12} = \frac{1}{12} \cdot \frac{1}{12} $	summation part of Eq. (25):	
user fälture sner acceptsner	Pe = Prof. Las		, 2 q ₁ p ₁ (1-p ₁) =t+1	see Eq. (15)	see Eq. (16)	see Eq. (25)	n Z PRAKORO
user faiture	PriE	±1.		$\sum_{j=1}^{N} q_{j-j-1} = \prod_{i=1}^{M} (1-\sigma_{ji}\beta_{1})!$	$M = (1 - \beta)^{M}$	$\frac{1}{1-1} q_j (1-(1-\sigma_j \beta_j))$:
Probability of	Model	Equal failure probabilities $P_1 = \frac{W}{N}$.	Unequal failure probabilities M=N, σ_{μ} = δ_{μ}	With statistical dependence	Special case of above	Random testing	Approximation of above

Types of Error Probability

MISSION of Rome Air Development Center

RANC plans and executes research, development, test and selected acquisition programs in support of Command, Control Communications and Intelligence (C^3I) activities. Technical and engineering support within areas of technical competence is provided to ESP Program Offices (POs) and other ESD elements. The principal technical mission areas are communications, electromagnetic guidance and control, surveillance of ground and aerospace objects, intelligence data collection and handling, information system technology, ionospheric propagation, solid state sciences, microwave physics and electronic reliability, maintainability and compatibility.